

Sort well with energy-constrained comparisons

Barbara Geissmann and Paolo Penna

Department of Computer Science, ETH Zurich, Zurich, Switzerland

Abstract

We study very simple sorting algorithms based on a *probabilistic* comparator model. In our model, errors in comparing two elements are due to (1) the energy or effort put in the comparison and (2) the difference between the compared elements. Such algorithms keep comparing pairs of randomly chosen elements, and they correspond to Markovian processes. The study of these Markov chains reveals an interesting phenomenon. Namely, in several cases, the algorithm which repeatedly compares only *adjacent* elements is *better* than the one making arbitrary comparisons: on the long-run, the former algorithm produces sequences that are “better sorted”. The analysis of the underlying Markov chain poses new interesting questions as the latter algorithm yields a non-reversible chain and therefore its stationary distribution seems difficult to calculate explicitly.

Digital Object Identifier [10.4230/LIPIcs...](https://doi.org/10.4230/LIPIcs...)

1 Introduction

Suppose one has to sort a number of elements by making pairwise comparisons, but sometimes the result of a comparison is *incorrect*. Sometimes the errors are unavoidable, and sometimes they are deliberately introduced in order to save other important resources. For example, with *probabilistic CMOS* it is possible to *trade energy for errors*, that is, one can reduce the energy spent for a single operation but this will increase the probability of incorrect response (see the survey of Palem and Lingamneni [17]). Errors also occur in measurements that require high precision (where a small noise can affect the result), or judgment made by individuals (who naturally tend to make small mistakes). One can envision the following situation:

1. It is “easier” to compare two elements if they differ “a lot”, while errors are more likely when they are “very close”.
2. If we are able or willing to spend more energy (effort) on a single comparison, then we can increase the probability of getting the correct result.

Given this scenario, one would like to design a strategy (algorithm) which sorts the elements nearly correctly. The following is thus a natural question:

Which strategies (algorithms) perform better?

This question has been already addressed under various models of errors (see e.g. [1, 6, 12]). The purpose of this work is to study this question under a new model which captures the two features above (Items 1 and 2).

1.1 Our contribution

In this work we propose to look at extremely simple *sorting* algorithms on what we call a *probabilistic comparator* model. These algorithms can be studied through the lens of *Markov chains* whose analysis reveals interesting properties regarding their behavior.



licensed under Creative Commons License CC-BY



Leibniz International Proceedings in Informatics

Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

Probabilistic comparator model

We introduce a simple model in which the probability that a comparison between two elements (numbers) is correct depends on two factors: (1) how different the two elements are and (2) the “effort” or “energy” spent to make the comparison. Intuitively speaking, the more energy the more accurate is the comparison, meaning that even very similar elements (small difference) can be distinguished. More precisely, for an energy parameter $\lambda \geq 1$, a comparison between two numbers a and b is correct with probability

$$p_{ab} := \frac{\lambda^{b-a}}{\lambda^{a-b} + \lambda^{b-a}},$$

where b is the biggest between a and b . One should think about a single comparison as a measurement of two quantities, which sometimes can be erroneous especially when the difference is small. All comparisons (including those involving the same two numbers) are independent and performed with the same parameter λ . That is, we consider *independent* errors in which the error probabilities are described by the formula above.

► **Remark.** Our model is inspired by the classical models in *statistical physics* [16] and in *game theory* [4] where $\lambda = e^{1/\text{noise}}$. In the latter applications, $1/\text{noise}$ corresponds to the “rationality level” of players, i.e., their ability to distinguish between strategies with similar payoffs. In statistical physics, the *noise* parameter is the “temperature” of the system, and high temperature corresponds to highly disordered configurations. To some extent, this model can be seen as an abstraction of the *probabilistic CMOS* technology which allows to trade energy for correctness [17]. As our model attempts to abstract from hardware-specific construction details, it necessarily introduces certain simplifying assumptions. Arguably, the major of these assumptions is that the probability of errors depends uniquely on the difference between the two numbers to be compared, and not on their actual values. On the other hand, the parameter λ captures the property that, by increasing the energy per single operation, the probability of correct comparisons increases in some way (depending on the hardware).

Sort by random comparisons

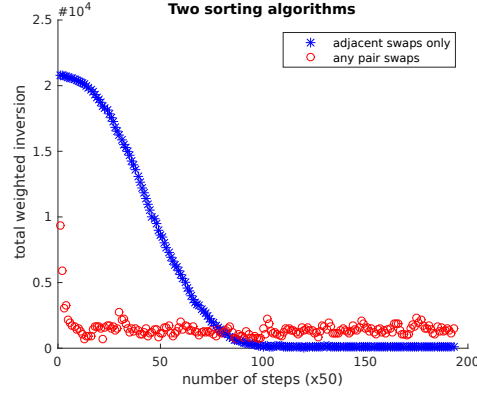
When sorting several elements, one performs several comparisons for a certain number of steps and then returns the resulting sequence. Consider the following two simple algorithms:

- *Any pairs swaps.* Compare two randomly chosen elements of the sequence.
- *Adjacent swaps.* Compare a randomly chosen element with the next one in the sequence.

Surprisingly enough, Figure 1 shows that the algorithm with only adjacent swaps gives *better* results after a certain number of comparisons are made (intuitively, the y -axis measures the “disorder” in the current sequence in each algorithm). Note that the initial input sequence is irrelevant as long as we consider sufficiently many steps of such algorithms, a property that can be formally captured by viewing these algorithms as Markov chains.

Algorithms and Markov chains

The two algorithms above (and others) can be viewed as Markov chains whose stationary distribution describes the output if we let them run long enough. We believe that these processes are interesting by themselves (for instance, they can be seen as variants of other well-studied processes – see Section 1.2), and they pose new questions on how to analyze them. Specifically, the adjacent swaps algorithm corresponds to a *reversible* Markov chain



■ **Figure 1** Comparison of two simple algorithms for sorting: on the long run the algorithm doing only adjacent comparisons gives a better result compared to the one doing all comparisons. The input sequence is $\{50, 49, \dots, 1\}$ while the sorted sequence is $\{1, 2, \dots, 50\}$. In this experiment we set $\lambda = e^{\frac{1}{5}}$, though the same behavior occurs for essentially any fixed λ (see Section 5 for the experiments).

\mathcal{M}_{adj} and its stationary distribution has a simple closed formula of the form

$$\pi(s) \propto \lambda^{-2w(s)} \quad (1)$$

where $w()$ is what we call a *total weighted inversion* of sequence s ,

$$w(s) := \sum_{i < j: s_i > s_j} s_i - s_j,$$

a measure of its “distance” from the correctly sorted sequence. Intuitively, inversions involving very different elements count more than inversions of almost identical elements. In contrast, the algorithm with arbitrary (random) pair comparisons corresponds to a *non-reversible* chain \mathcal{M}_{any} , and therefore the analysis of its stationary distribution is considerably more complicated. We also provide a variant of \mathcal{M}_{any} in which comparisons of non-adjacent pairs are done *multiple times*: For example, if numbers a and b are two positions away in the current sequence (say $a = s_i$ and $b = s_{i+2}$) then we perform *two* comparisons and accept to swap them only if *both* of them tell to do so. This third chain \mathcal{M}_{any}^* has the same stationary distribution of \mathcal{M}_{adj} . Therefore, one can see this “careful swapping” rule as a way to fix the algorithm doing naively arbitrary swaps. The analysis of this chain \mathcal{M}_{any}^* is based on the *Kolmogorov reversibility criterion*.

Figure 1, and all experiments we made on various input sequences, suggest that \mathcal{M}_{adj} yields *better sorted* sequences than \mathcal{M}_{any} , though the latter chain *converges faster* to its stationary distribution. We study both the *mixing time* and the properties of the *stationary distribution*, like the probability of returning the sorted sequence.

- In Section 3, we consider the case of *binary sequences*, where each element in the sequence is either a or b . We show that the mixing time of \mathcal{M}_{adj} is $O(n^2)$, while for \mathcal{M}_{any} it is $O(n \log n)$, or even linear if the number of occurrences of b is constant, for every $\lambda > 1$.
- In Section 4 we study the probability that the chains return the sorted sequence at stationary distribution. We show that \mathcal{M}_{adj} is better than \mathcal{M}_{any} when sorting *three arbitrary elements*. This result is based on the *Markov chain tree theorem* and it is the most involved in this section. Similar results hold also for sorting arbitrary long sequences with a *single outlier*, that is, binary sequences with a single element $b > a$ and

many a 's. Here the analysis shows a quantitative difference between the two chains (cf. Theorem 14). Note that, all these results apply also to \mathcal{M}_{any}^* in place of as \mathcal{M}_{adj} , since they have the same stationary distribution.

1.2 Related work

Stochastic models of the form (1) are very common in statistics and, in particular, Mallows [15] was among the firsts to consider such models in the context of permutations: There the weight function $w()$ is a suitable distance function which comes from probabilities p_{ab} of ranking a before b . In that sense, our model is a special case of Mallows' model, though the procedure of [15] is different: One makes all pairwise comparisons at once until a consistent result is obtained. Our probabilistic comparator is also a special case of Bradley and Terry [5] where the probability p_{ab} of ranking a before b is of the form $\frac{w_a}{w_a + w_b}$.

Several restrictions on p_{ab} have been studied for the natural Markov chain which makes only adjacent comparisons. The classical card shuffling problem corresponds to the *unbiased* version of this chain in which all probabilities p_{ab} equal $1/2$, for which Wilson [18] proved that this chain is rapidly mixing and gave a very tight bound. A similar problem is the uniform sampling of partial order extensions, which corresponds to probabilities p_{ab} being $1/2$ or 1 and $p_{ba} = 1 - p_{ab}$. For the latter, Bubley and Dyer [7] showed that this chain is also rapidly mixing. Benjamini et al. [2] proved rapidly mixing for the *constant biased* case, that is, when every comparison is correct with some fixed probability $p > 1/2$, independently of the compared elements: $p_{ab} = p > 1/2$ for all $a < b$. The mixing time of *biased* comparisons has been studied by Bhakta et al. [3] under two comparison models called “choose your weapon” and “league hierarchies”: In the first model p_{ab} depends only on the largest between a and b , while in the second model all numbers are the leaves of some tree and p_{ab} depends only on the least common ancestor of a and b . Note that our model does not fall in either class even for the case of only three distinct elements.

Diaconis and Ram [8] studied a different type of chains called systematic scan algorithms: for the *unbiased* case, they proved that n of such scan operations are sufficient to reach the stationary distribution.

1.3 Preliminary definitions on Markov chains

In this section we introduce some of the definitions on Markov chains used throughout this work (for more details see Levin et al. [14]). A Markov chain over a finite state space S is specified by a transition matrix P , where $P(s, s')$ is the probability of moving from state s to state s' in one step. The t^{th} power of the transition matrix gives the probability of moving from one state to another state in t steps. All chains studied in this work are ergodic meaning that they have a unique *stationary distribution* π : for any two states s and s' , $\lim_{t \rightarrow \infty} P^t(s, s') = \pi(s')$.

We will use the definition of a *reversible* Markov chain, also called *detailed balanced condition*: If the transition matrix P admits a vector π such that $\pi(s)P(s, s') = \pi(s')P(s', s)$ for all s and s' , then π is the *stationary distribution* of the chain with transitions P .

An equivalent characterization of reversible chains is given by looking at cycles over the states. For any subset $\Gamma \subseteq S \times S$ of transitions (pairs of states of the chain), define the associated probability as the product of all these transitions in the chain,

$$\mathbf{P}(\Gamma) := \prod_{(x,y) \in \Gamma} P(x,y) . \quad (2)$$

Let Γ^{-1} denote the reversed edges in Γ , that is, $\Gamma^{-1} := \{(y, x) \mid (x, y) \in \Gamma\}$. The *Kolmogorov reversibility criterion* [13] says that a chain is reversible if and only if for any cycle C ,

$$\mathbf{P}(C) = \mathbf{P}(C^{-1}) . \quad (3)$$

For the sake of clarity, we sometimes denote cycles as $C = s^1 \rightarrow s^2 \rightarrow \dots \rightarrow s^\ell \rightarrow s^1$ and the corresponding reversal by $C^{-1} = s^1 \leftarrow s^2 \leftarrow \dots \leftarrow s^\ell \leftarrow s^1$.

The stationary distribution of any (even non-reversible) Markov chain can be computed by looking at the probabilities of all directed trees rooted at some state. More formally, let $\mathcal{T}(s)$ be the set of all directed trees rooted at state s , that is, from every other state there is a path towards s in the tree. The *Markov chain tree theorem* (see Freidlin and Wentzell [10], Chapter 6, Lemma 3.1) says that, for any ergodic Markov chain with transition matrix P , its stationary distribution π is given by:

$$\pi(s) = \frac{W(s)}{\sum_{\hat{s}} W(\hat{s})} , \quad \text{where} \quad W(s) := \sum_{T \in \mathcal{T}(s)} \mathbf{P}(T) . \quad (4)$$

1.4 Measure of Disorder

In this section we introduce a formal definition for the *total weighted inversion*, which can be seen as a measure of disorder. As we shall prove in the next section, this arises naturally from the algorithm performing adjacent swaps.

► **Definition 1.** The *total weighted inversion* of a sequence s is defined as

$$w(s) := \sum_{i < j: s_i > s_j} s_i - s_j .$$

► **Example 2.** Consider the sequence $s = (5, 2, 3)$ and the sorted sequence $(2, 3, 5)$. Then the total weighted inversion of s is equal to $w(s) = (5 - 2) + (5 - 3) = 5$.

The displacements of the single elements allow an equivalent way to describe the total weighted inversion (this equivalent definition turns out to be useful in the next section).

► **Lemma 3.** For a sequence s let $s^{(sort)}$ be the sequence sorted in non-decreasing order. Then, $w(s) = \sum_i (s_i^{(sort)} - s_i)i$.

Proof. In the sum $\sum_{i < j: s_i > s_j} s_i - s_j$, every element s_i is added r_i and subtracted l_i times, where r_i is the number of smaller elements on its right hand side and l_i the number of larger elements on its left. The difference $d_i = r_i - l_i$ corresponds to the displacement of s_i to the right compared to the sorted sequence, i.e., $s_{i+d_i}^{(sort)} = s_i$. Since $d_i \cdot s_i$ is exactly the contribution of s_i to $w(s)$, the claim follows immediately. ◀

2 Sorting algorithms as Markov chains

In this section we define the algorithms and the resulting Markov chains. The first chain performs only adjacent comparisons.

► **Definition 4.** The chain \mathcal{M}_{adj} is defined as follows:

1. Pick an index i in $\{1, \dots, n-1\}$ uniformly at random;
2. Swap $s_i = a$ and $s_{i+1} = b$ with probability $\frac{\lambda^{a-b}}{\lambda^{a-b} + \lambda^{b-a}}$.

XX:6 Sort well with energy-constrained comparisons

We shall prove below that the stationary distribution of this chain assigns higher probabilities to the sequences that are “nearly sorted”.

► **Theorem 5.** *The chain \mathcal{M}_{adj} is reversible with stationary distribution $\pi(s) \propto \lambda^{-2w(s)}$, where $w(s)$ is the total weighted inversion.*

Proof. We prove that the chain is reversible. Let s and s' be two sequences that differ in i 's swap (otherwise $P(s, s') = 0 = P(s', s)$ and reversibility is trivial). Observe that by definition

$$\frac{P(s, s')}{P(s', s)} = \lambda^{2(s_i - s_{i+1})} = \lambda^{2(a-b)} \quad \text{and} \quad \frac{\pi(s')}{\pi(s)} = \lambda^{2w(s) - 2w(s')} .$$

Since $s^{(sort)} = s'^{(sort)}$ and s' is obtained by swapping $a = s_i$ and $b = s_{i+1}$,

$$\begin{aligned} w(s) - w(s') &= (s_i^{(sort)} - s_i)i - (s_i^{(sort)} - s'_i)i \\ &\quad + (s_{i+1}^{(sort)} - s_{i+1})(i+1) - (s_{i+1}^{(sort)} - s'_{i+1})(i+1) = a - b \end{aligned}$$

and therefore the detailed balance condition is satisfied. ◀

We next consider chains which compare any two numbers in the sequence:

► **Definition 6.** The chain \mathcal{M}_{any} is defined as follows:

1. Pick two indexes i and j in $\{1, \dots, n\}$ uniformly at random, with $i < j$;
2. Swap $s_i = a$ and $s_j = b$ with probability $\frac{\lambda^{a-b}}{\lambda^{a-b} + \lambda^{b-a}}$.

The chain \mathcal{M}_{any}^* is defined as above except that the probability of swapping is $\left(\frac{\lambda^{a-b}}{\lambda^{a-b} + \lambda^{b-a}}\right)^{j-i}$.

► **Theorem 7.** *Chain \mathcal{M}_{any}^* is reversible and has the same stationary distribution as \mathcal{M}_{adj} , while chain \mathcal{M}_{any} is not reversible.*

Proof. To prove that \mathcal{M}_{any}^* has the same stationary distribution as \mathcal{M}_{adj} we argue as follows. Consider any transition from s to s' which swaps two elements at distance $k \geq 1$. There exists a path \mathcal{P} in \mathcal{M}_{adj} that leads from s to s' and whose probability has the same form of the single transition in \mathcal{M}_{any}^* ,

$$\mathbf{P}(\mathcal{P}) = \prod_{(x,y) \in \mathcal{P}} P(x, y) = \frac{\lambda^{k(a-b)}}{D(\mathcal{P})} .$$

The path is obtained by simulating the swap between a and b via adjacent swaps:

$$\begin{aligned} \mathcal{P} := & (\dots a x_1 \dots x_{k-1} b \dots) \leftarrow (\dots x_1 a \dots x_{k-1} b \dots) \dots \leftarrow (\dots x_1 \dots x_{k-1} a b \dots) \leftarrow \\ & (\dots x_1 \dots x_{k-1} b a \dots) \leftarrow (\dots x_1 x_2 \dots b x_{k-1} a \dots) \dots \leftarrow (\dots b x_1 x_2 \dots x_{k-1} a \dots) \end{aligned}$$

which yields in the numerator the product

$$\lambda^{(a-x_1) + \dots + (a-x_{k-1}) + (a-b) + (x_{k-1}-b) + \dots + (x_1-b)} = \lambda^{k(a-b)} .$$

Note also that the reverse path \mathcal{P}^{-1} leading from s' to s has probability $\mathbf{P}(\mathcal{P}^{-1}) = \frac{\lambda^{k(b-a)}}{D(\mathcal{P})}$, where the denominator $D(\mathcal{P})$ is the same as above because all transitions in the chains are of the form $P(x, y) = N_{xy}/D_{xy}$ with $D_{xy} = D_{yx}$. Since \mathcal{M}_{adj} is reversible, we get the first of the following equalities:

$$\frac{\pi(s')}{\pi(s)} = \frac{\mathbf{P}(\mathcal{P})}{\mathbf{P}(\mathcal{P}^{-1})} = \frac{\lambda^{k(a-b)}}{\lambda^{k(b-a)}} = \frac{P^*(s, s')}{P^*(s', s)} ,$$

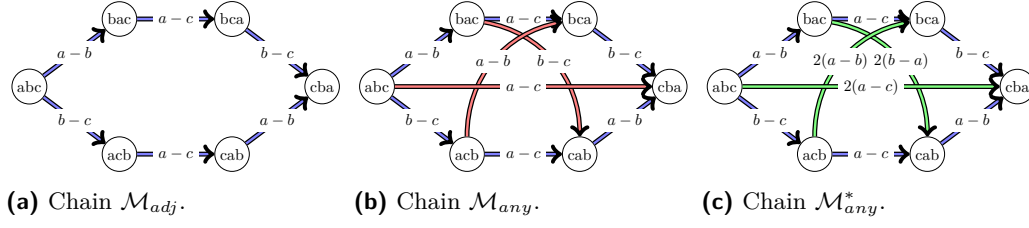


Figure 2 The three chains for sorting three elements abc ; A transition with label w has probability $\frac{\lambda^w}{\lambda^w + \lambda - w}$; For clarity sake we only show forward transitions.

where P^* is the transition matrix of \mathcal{M}_{any}^* . Thus, the detailed balance condition for P^* is satisfied and π is the stationary distribution of \mathcal{M}_{any}^* .

Finally, to see that \mathcal{M}_{any} is not reversible, consider the cycle $(abc) \leftarrow (bac) \leftarrow (bca) \leftarrow (cba) \leftarrow (cab) \leftarrow (acb) \leftarrow (abc)$ in Figure 2b which has a length different from the reversed cycle. This violates the Kolmogorov reversibility criterion (3). \blacktriangleleft

In our experiments (see Section 5), it turns out that doing only adjacent comparisons (\mathcal{M}_{adj}) is better than doing any comparisons (\mathcal{M}_{any}); The following sections provide analytical results for special cases. Note that Theorem 7 says that in the long-run \mathcal{M}_{any}^* is as good as \mathcal{M}_{adj} .

3 Binary inputs

In this section we restrict to the case in which every element in the sequence is either a or b for some $b > a$. That is, the sorted sequence is $(a, \dots, a, b, \dots, b)$, where n_a denotes the number of a 's and n_b denotes the number of b 's.

3.1 Mixing time

For binary inputs, we can uniquely express every sequence by a vector $v \in \{0, \dots, n_a\}^{n_b}$, $v_1 \geq v_2 \geq \dots \geq v_{n_b}$, where v_i denotes the number of inversions of the i -th b in the sequence (for example, $babba$ corresponds to 211, while $babab$ corresponds to 210). Such a vector is visualized as a monotonically decreasing ‘staircase’ in a $n_b \times n_a$ grid and \mathcal{M}_{adj} corresponds to the biased Markov process in [11]. The bounds on the mixing time for this process translate immediately for our chain.

► **Theorem 8** (by Theorem 2.1 in [11]). *For binary inputs, the mixing time of \mathcal{M}_{adj} satisfies*

$$t_{mix}(\epsilon) = O(n^2 \log(\epsilon^{-1})).$$

Observe also that the chain \mathcal{M}_{adj} corresponds to the well-known *asymmetric simple exclusion process* (see [2] and [3]).

We next consider \mathcal{M}_{any} and prove an upper bound. To bound the mixing time of \mathcal{M}_{any} we use the method of *path coupling* [9]. A *path coupling* for a chain \mathcal{M} can be specified by providing distributions

$$\mathbb{P}_{x,y}[X = x', Y = y'], \quad \text{for all } x, y \in S \text{ such that } P(x, y) > 0, \quad (5)$$

satisfying, for all $x, y \in S$ such that $P(x, y) > 0$,

$$\mathbb{P}_{x,y}[X = x'] = P(x, x') \quad \text{for all } x' \in S, \quad (6)$$

$$\mathbb{P}_{x,y}[Y = y'] = P(y, y') \quad \text{for all } y' \in S. \quad (7)$$

XX:8 Sort well with energy-constrained comparisons

We use ρ to denote the shortest-path distance in the Markov chain, i.e., $\rho(x, y)$ is the minimum number of transitions to go from x to y .

► **Lemma 9** (Theorem 2.1 in [9]). *Suppose there exists $\beta < 1$ such that, for all x, y with $P(x, y) > 0$, it holds that*

$$\mathbb{E}_{x,y}[\rho(X, Y)] \leq \beta. \quad (8)$$

Then the mixing time $t_{\text{mix}}(\epsilon)$ of the Markov chain under consideration satisfies

$$t_{\text{mix}}(\epsilon) \leq \frac{\log(D\epsilon^{-1})}{1 - \beta}.$$

Path coupling for \mathcal{M}_{any} .

Consider two sequences x and y which differ by swapping elements in position i^* and j^* . For every such pair (x, y) we specify the probabilities in (5) to move to a pair (x', y') . We group the $\binom{n}{2}$ different swaps between elements in positions i and j as follows:

$$(i^*, j^*) \leftrightarrow (i^*, j^*) \quad (i^*, j) \leftrightarrow (j^*, j) \quad (i, j^*) \leftrightarrow (i, i^*) \quad (i, j) \leftrightarrow (i, j),$$

in the sense that if we consider the positions i^* and j in one sequence, we consider the positions j^* and j in the other sequence, and vice versa. Clearly, this defines a bijection on the swaps of the two sequences. Now let $x_{i \times j}$ denote the sequence obtained from x by swapping the two elements at positions i and j . The path coupling is as follows:

for i^*, j^*	$(x, y) \mapsto (y, y)$	with $P(x, y)$,
	$(x, y) \mapsto (x, x)$	with $P(y, x)$,
for i^*, j	$(x, y) \mapsto (x_{i^* \times j}, y_{j^* \times j})$	with $\min\{P(x, x_{i^* \times j}), P(y, y_{j^* \times j})\}$,
	$(x, y) \mapsto (x_{i^* \times j}, y)$	with $\max\{0, P(x, x_{i^* \times j}) - P(y, y_{j^* \times j})\}$,
	$(x, y) \mapsto (x, y_{j^* \times j})$	with $\max\{0, P(y, y_{j^* \times j}) - P(x, x_{i^* \times j})\}$,
for i, j^*	$(x, y) \mapsto (x_{i \times j^*}, y_{i \times i^*})$	with $\min\{P(x, x_{i \times j^*}), P(y, y_{i \times i^*})\}$,
	$(x, y) \mapsto (x_{i \times j^*}, y)$	with $\max\{0, P(x, x_{i \times j^*}) - P(y, y_{i \times i^*})\}$,
	$(x, y) \mapsto (x, y_{i \times i^*})$	with $\max\{0, P(y, y_{i \times i^*}) - P(x, x_{i \times j^*})\}$,
for i, j	$(x, y) \mapsto (x_{i \times j}, y_{i \times j})$	with $P(x, x_{i \times j}) = P(y, y_{i \times j})$.

Finally, with all remaining probability

$$(x, y) \mapsto (x, y).$$

One can easily check that this is indeed a path coupling, that is, (6)-(7) are satisfied. The difficulty is in proving the condition necessary to apply Lemma 9.

► **Lemma 10.** *The path coupling defined above satisfies condition (8) with*

$$\beta \leq 1 - \frac{2(1 + p(n-2))}{n(n-1)} \leq 1 - \frac{2p}{n-1}.$$

Proof. The second inequality follows from $p(n-2) + 1 \geq pn$, since $p \leq \frac{1}{2}$. We next prove the first inequality. Let x, y be two sequences that differ in swapping positions i^* and j^* , thus $\rho(x, y) = 1$. Since $P(x, y) + P(y, x) = 1$, the new distance $\rho(x', y')$ after choosing

positions i^* and j^* is always zero. Furthermore, for every position k , such that $k \neq i^*$ and $k \neq j^*$, either $\rho(x_{i^* \times k}, y_{j^* \times k}) = 0$ or $\rho(x_{k \times j^*}, y_{k \times i^*}) = 0$, and the probability of accepting such a transition is at least p . Finally, it is easy to see that after every other transition $\rho(x', y') = 1$. Remember that there are $\binom{n}{2}$ pairs of positions in a sequence. Therefore,

$$\mathbb{E}[\rho(x', y')] \leq \left(1 - \frac{1+p(n-2)}{\binom{n}{2}}\right). \quad \blacktriangleleft$$

► **Theorem 11.** *Let n_a (resp., n_b) denote the number of a 's (resp., b 's) in the sequence. The mixing time of \mathcal{M}_{any} satisfies*

$$t_{mix}(\epsilon) \leq \frac{n(\log(n') - \log(\epsilon))}{2p},$$

where $n' = \min\{n_a, n_b\} \leq \frac{n}{2}$.

Proof. The diameter D is the maximum number of transitions required to go from any sequence to any other sequence. Then $D = \min\{n_a, n_b\} \leq \frac{n}{2}$, because with D swaps we can either move all a or all b to their desired positions. By Lemmata 9 and 10, the claim follows immediately. \blacktriangleleft

4 Only adjacent swaps is better

In this section we prove that, for two special cases, the chain \mathcal{M}_{adj} performing only adjacent comparisons is better than the chain \mathcal{M}_{any} performing comparisons between any two elements.

4.1 Three elements

Our first special case is to consider sorting three arbitrary elements and show that \mathcal{M}_{adj} has more chances to return the sorted sequence than \mathcal{M}_{any} .

► **Theorem 12.** *For any three elements, not all of them identical, the chain \mathcal{M}_{adj} returns the sorted sequence with a probability (at stationary distribution) strictly larger than that of \mathcal{M}_{any} (at stationary distribution),*

$$\pi_{adj}(abc) > \pi_{any}(abc) \quad (9)$$

where abc is the sorted sequence ($a \leq b \leq c$), for all $\lambda > 0$.

In order to prove this theorem we show that the ratios between the distribution of adjacent states in \mathcal{M}_{adj} gets “worse” in \mathcal{M}_{any} :

► **Lemma 13.** *For any two states s and s' that differ in exactly one adjacent swap, if the total weighted inversion satisfies $w(s') > w(s)$, then it holds that*

$$\frac{\pi_{adj}(s')}{\pi_{adj}(s)} < \frac{\pi_{any}(s')}{\pi_{any}(s)}. \quad (10)$$

Proof Idea. We use the Markov Chain Tree Theorem (4). Our goal is thus to show that

$$\frac{\sum_{T \in \mathcal{T}(s)} \mathbf{P}(T)}{\sum_{T' \in \mathcal{T}(s')} \mathbf{P}(T')} < \frac{\pi_{adj}(s)}{\pi_{adj}(s')} = \lambda^{2(w(s') - w(s))} \quad \text{for} \quad w(s) < w(s'). \quad (11)$$

Ideally, one would like to find a bijection from trees $T \in \mathcal{T}(s)$ to trees $T' \in \mathcal{T}(s')$ such that $\frac{\mathbf{P}(T)}{\mathbf{P}(T')} < \frac{\pi_{adj}(s)}{\pi_{adj}(s')}$ holds for each tree $T \in \mathcal{T}(s)$. Unfortunately, this is in general not possible, so the following slightly more involved argument is used:

XX:10 Sort well with energy-constrained comparisons

- The simple mapping we use consists in reversing the path from s' to s in T to obtain the new tree T' . This mapping is a bijection between $\mathcal{T}(s)$ and $\mathcal{T}(s')$.
- Because this mapping does not guarantee the desired inequality for all trees T , we classify the trees in $\mathcal{T}(s)$ into *good* and *bad* trees: a tree T is *bad* if $\frac{\mathbf{P}(T)}{\mathbf{P}(T')} > \lambda^{2(w(s')-w(s))}$ and *good* otherwise. We then show that

$$\frac{\sum_{T \in \text{bad}} \mathbf{P}(T) + \sum_{T \in \text{good}} \mathbf{P}(T)}{\sum_{T \in \text{bad}} \mathbf{P}(T') + \sum_{T \in \text{good}} \mathbf{P}(T')} < \lambda^{2(w(s')-w(s))} ,$$

where T' is the tree obtained from T via the mapping in the previous item. This proves (11) since *good* and *bad* define a partition of $\mathcal{T}(s)$ and also a partition of $\mathcal{T}(s')$.

The details of this proof are given in Appendix B.1. ◀

From this lemma it is easy to obtain (9) in Theorem 12.

Proof of Theorem 12. By transitivity, Lemma 13 implies that, for all non-sorted sequences $s \neq (abc)$, $\frac{\pi_{adj}(s)}{\pi_{adj}(abc)} < \frac{\pi_{any}(s)}{\pi_{any}(abc)}$, and therefore

$$\pi_{adj}(abc) = \frac{1}{1 + \sum_{s \neq (abc)} \frac{\pi_{adj}(s)}{\pi_{adj}(abc)}} > \frac{1}{1 + \sum_{s \neq (abc)} \frac{\pi_{any}(s)}{\pi_{any}(abc)}} = \pi_{any}(abc).$$

◀

4.2 One outlier

We call *one outlier* the case in which we have $n - 1$ small identical elements, and only one bigger element (the outlier) to be sorted. That is, the sorted sequence is (a, a, \dots, a, b) , with $b > a$. Since swapping two identical elements does not change the sequence, i.e. the state of the chain, we have n states which correspond to the possible positions of b . We denote the state in which b is in position i by $s^{(i)}$, so that $s^{(n)} = (a, \dots, a, b)$ is the sorted sequence and $s^{(1)} := (b, a, \dots, a)$ is the “reversely sorted” sequence. Note that the weighted inversion of $s^{(i)}$ is $(n - i)(b - a)$, thus implying that the expected total weighted inversion is of the form

$$E^w := \sum_{i=1}^n (n - i)(b - a)\pi(s^{(i)}) . \quad (12)$$

It is useful for the analysis to consider the probability that element b is erroneously declared smaller than a in a single comparison,

$$p := 1 - p_{ab} = \frac{\lambda^{a-b}}{\lambda^{a-b} + \lambda^{b-a}} .$$

► **Observation 1.** For the one outlier case, the chain \mathcal{M}_{adj} becomes a so-called birth-and-death chain meaning that from each state $s^{(i)}$ we can only move to $s^{(i+1)}$ or to $s^{(i-1)}$, and the transition probabilities are

$$P(s^{(i+1)}, s^{(i)}) = \frac{p}{n-1} , \quad P(s^{(i)}, s^{(i+1)}) = \frac{1-p}{n-1} ,$$

for all $i \in \{1, \dots, n-1\}$. In the chain \mathcal{M}_{any} every state is connected to all other states and the transition probabilities are $P(s^{(i)}, s^{(j)}) = p/\binom{n}{2}$ if $i > j$ and $(1-p)/\binom{n}{2}$ if $i < j$.

The next theorem says that the chain \mathcal{M}_{adj} has a better probability of returning the sorted sequence and a better expected total weighted inversion than \mathcal{M}_{any} .

► **Theorem 14.** *For the case of one outlier, the following holds. The probability of obtaining the sorted sequence, at stationary distribution, is constant for the \mathcal{M}_{adj} , while for \mathcal{M}_{any} it converges to zero as n grows:*

$$\pi_{adj}(s^{(sorted)}) > \frac{1-2p}{1-p} \quad , \quad \pi_{any}(s^{(sorted)}) < \frac{1-p}{np} \quad .$$

The expected total weighted inversion is constant for \mathcal{M}_{adj} , while for \mathcal{M}_{any} it grows linearly in n :

$$E_{adj}^w < \frac{p}{1-2p}(b-a) \quad , \quad E_{any}^w > np(b-a) \quad .$$

The theorem above follows immediately from the next two lemmas, whose proof is given in Appendix B.2. Recall that $s^{(sorted)} = s^{(n)}$.

► **Lemma 15.** *The stationary distributions of \mathcal{M}_{adj} and \mathcal{M}_{any} , are*

$$\pi_{adj}(s^{(i)}) = \frac{p^{n-i}(1-p)^{i-1}(1-2p)}{(1-p)^n - p^n} \quad ,$$

$$\pi_{any}(s^{(i)}) = \frac{np(1-p)}{((n-i+1)(1-p) + (i-1)p)((n-i)(1-p) + ip)} \quad .$$

► **Lemma 16.** *For \mathcal{M}_{adj} and \mathcal{M}_{any} , the corresponding expected total weighted inversions are*

$$E_{adj}^w = n(b-a)p \left(\frac{1}{n(1-2p)} - \frac{p^{n-1}}{(1-p)^n - p^n} \right) < (b-a) \frac{p}{1-2p} \quad ,$$

$$E_{any}^w = n(b-a)p \cdot \sum_{i=0}^{n-1} \frac{i(1-p)}{((i+1)(1-p) + (n-i-1)p)(i(1-p) + (n-i)p)} > n(b-a)p \quad .$$

5 Experimental results

We conducted the following set of experiments on several input sequences to compare the two sorting algorithms \mathcal{M}_{adj} and \mathcal{M}_{any} .

- **Low energy (high noise) regime.** We evaluate how much the two algorithms are robust to an increase of the error probability by taking $\lambda = e^{1/noise}$ for increasing values of *noise*. Figure 3 shows that \mathcal{M}_{any} degrades much earlier than \mathcal{M}_{adj} .
- **The best of the two.** We compare all three algorithms \mathcal{M}_{adj} , \mathcal{M}_{any} and \mathcal{M}_{any}^* in Figure 4. These experiments suggest that \mathcal{M}_{any}^* possesses good features from both the other algorithms: the total weighted inversion decreases faster than \mathcal{M}_{adj} , while its stationary distribution is of course better than \mathcal{M}_{any} , which is evident in the second range of noise.
- **Probability of getting sorted sequence.** We evaluate how the probability of hitting the sorted sequence changes when the noise increases. In Appendix A, Figure 5 deals with the sequence of ten elements $\{1, 2, \dots, 10\}$, while Figure 6 is about the one outlier $\{1, 1, 1, 1, 1, 1, 1, 1, 1, 2\}$.

XX:12 Sort well with energy-constrained comparisons

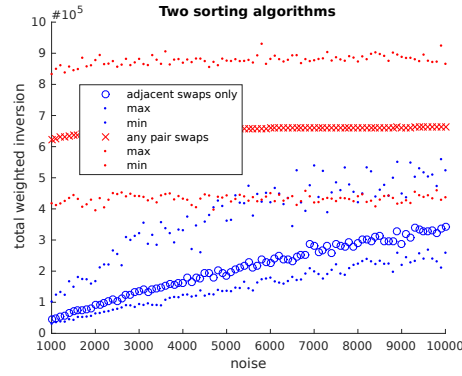
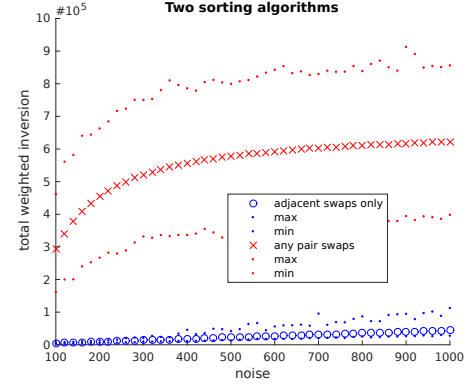
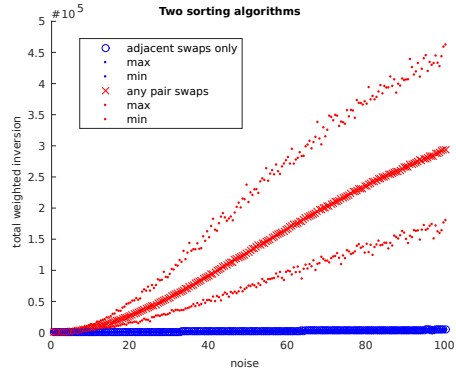


Figure 3 A comparison of \mathcal{M}_{adj} and \mathcal{M}_{any} on the long run. We measure the average, maximum, and minimum total weighted inversion for a certain number of steps after both algorithms have approached their stationary distribution. The elements to be sorted are $(200, 199, \dots, 1)$. By increasing the value of $noise$, where $\lambda = e^{1/noise}$, the stationary behavior of \mathcal{M}_{any} degrades much earlier than that of \mathcal{M}_{adj} .

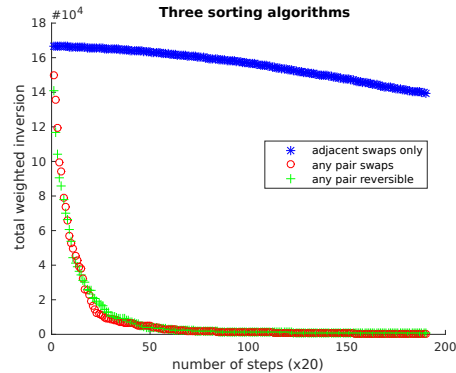
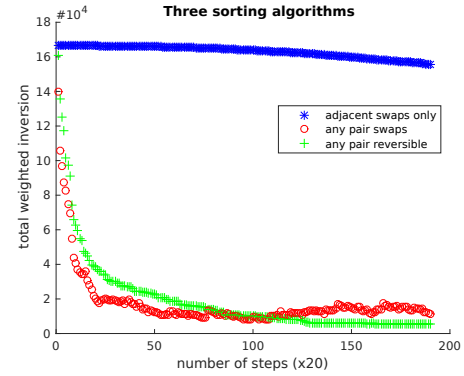
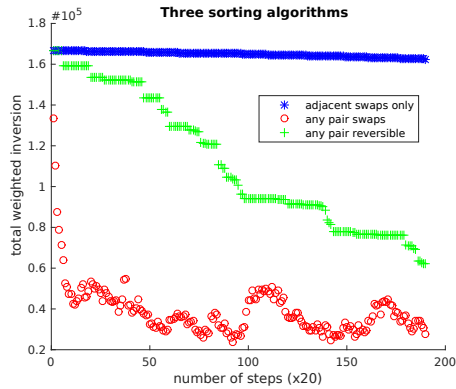


Figure 4 A comparison of the three algorithms \mathcal{M}_{adj} , \mathcal{M}_{any} , and \mathcal{M}_{any}^* . The elements to be sorted are $(100, 99, \dots, 1)$. The plots are for three different values of $noise$ $\{50, 20, 0.1\}$, where $\lambda = e^{1/noise}$.

Acknowledgments.

We are grateful to Lucas Boczkowski, Tomáš Gavenčíak, Francesco Pasquale, and Peter Widmayer for useful discussions. Research supported by SNF (project number 200021_165524).

References

1. Laurent Alonso, Philippe Chassaing, Florent Gillet, Svante Janson, Edward M Reingold, and René Schott. Quicksort with unreliable comparisons: a probabilistic analysis. *Combinatorics, Probability and Computing*, 13(4-5):419–449, 2004.
2. Itai Benjamini, Noam Berger, Christopher Hoffman, and Elchanan Mossel. Mixing times of the biased card shuffling and the asymmetric exclusion process. *Transactions of the American Mathematical Society*, 357(8):3013–3029, 2005.
3. Prateek Bhakta, Sarah Miracle, Dana Randall, and Amanda Pascoe Streib. Mixing times of Markov chains for self-organizing lists and biased permutations. In *Proc. of the 24th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 1–15, 2013.
4. Lawrence E. Blume. The statistical mechanics of strategic interaction. *Games and Economic Behavior*, 5:387–424, 1993.
5. Ralph Allan Bradley and Milton E. Terry. Rank Analysis of Incomplete Block Designs: I. The Method of Paired Comparisons. *Biometrika*, 39(3/4):324–345, 1952. ISSN 00063444. URL <http://www.jstor.org/stable/2334029>.
6. Mark Braverman and Elchanan Mossel. Noisy sorting without resampling. In *Proceedings of the nineteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 268–276. Society for Industrial and Applied Mathematics, 2008.
7. Russ Bubley and Martin Dyer. Faster random generation of linear extensions. *Discrete mathematics*, 201(1):81–88, 1999.
8. Persi Diaconis and Arun Ram. Analysis of systematic scan Metropolis algorithms using Iwahori-Hecke algebra techniques. *Michigan Math. J.*, 48(1):157–190, 2000.
9. Martin E. Dyer and Catherine S. Greenhill. A more rapidly mixing markov chain for graph colorings. *Random Struct. Algorithms*, 13(3-4):285–317, 1998.
10. M. Freidlin and A.D. Wentzell. *Random Perturbations of Dynamical Systems*. New York: Springer Verlag, 1984.
11. Sam Greenberg, Amanda Pascoe, and Dana Randall. Sampling biased lattice configurations using exponential metrics. In *Proceedings of the Twentieth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2009, New York, NY, USA, January 4-6, 2009*, pages 76–85, 2009.
12. Petros Hadjicostas and KB Lakshmanan. Recursive merge sort with erroneous comparisons. *Discrete Applied Mathematics*, 159(14):1398–1417, 2011.
13. Frank P. Kelly. *Reversibility and stochastic networks*. Cambridge University Press, 2011.
14. David Asher Levin, Yuval Peres, and Elizabeth Lee Wilmer. *Markov chains and mixing times*. American Mathematical Soc., 2009.
15. Colin L. Mallows. Non-null ranking models. I. *Biometrika*, 44(1/2):114–130, 1957.
16. Marc Mezard and Andrea Montanari. *Information, physics, and computation*. Oxford University Press, 2009.
17. Krishna V. Palem and Avinash Lingamneni. Ten Years of Building Broken Chips: The Physics and Engineering of Inexact Computing. *ACM Trans. Embedded Comput. Syst.*, 12(2s):87, 2013.
18. David Bruce Wilson. Mixing times of lozenge tiling and card shuffling Markov chains. *Annals of Applied Probability*, pages 274–325, 2004.

A Additional experiments

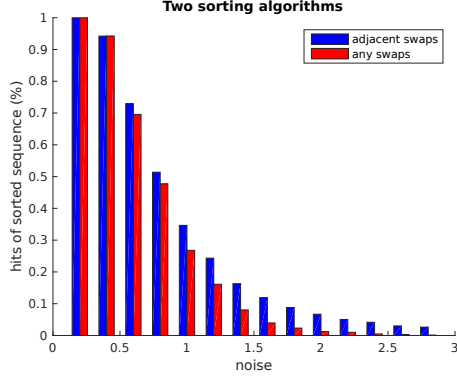


Figure 5 We measure the probability of hitting the sorted sequence during a certain number of steps after both algorithms have approached their stationary distribution. The elements to be sorted are $(10, 9, \dots, 1)$. By increasing the value of *noise*, where $\lambda = e^{1/\text{noise}}$, the stationary probability of sorted sequence decreases much faster in \mathcal{M}_{any} than in \mathcal{M}_{adj} .

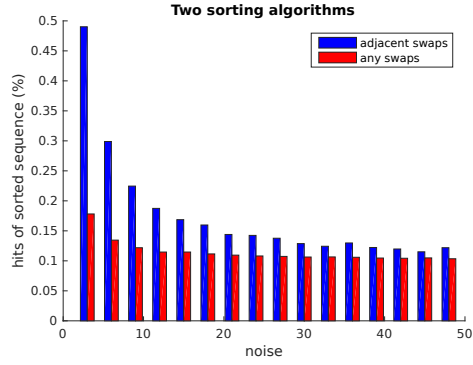
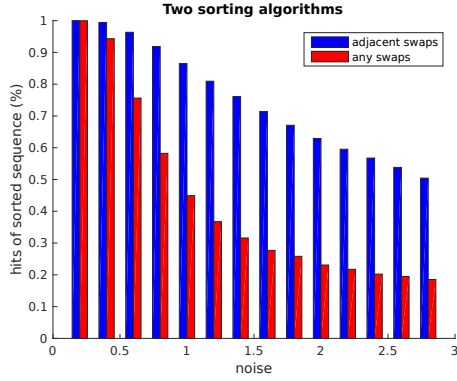


Figure 6 Percentage of hits of the sorted sequence after both algorithms have reached their stationary distribution. The set of elements is $\{1, 1, 1, 1, 1, 1, 1, 1, 2\}$.

B Additional lemmata and proofs

B.1 Three elements

► **Lemma 13.** *For any two states s and s' that differ in exactly one adjacent swap, if the total weighted inversion satisfies $w(s') > w(s)$, then it holds that*

$$\frac{\pi_{adj}(s')}{\pi_{adj}(s)} < \frac{\pi_{any}(s')}{\pi_{any}(s)} . \quad (10)$$

B.1.1 Proof of Lemma 13

We shall use the Markov Chain Tree Theorem (4). Our goal is thus to show that

$$\frac{\sum_{T \in \mathcal{T}(s)} \mathbf{P}(T)}{\sum_{T' \in \mathcal{T}(s')} \mathbf{P}(T')} < \frac{\pi_{adj}(s)}{\pi_{adj}(s')} = \lambda^{2(w(s') - w(s))} \quad \text{for} \quad w(s) < w(s') . \quad (13)$$

The strategy to prove this inequality is to find a suitable mapping from trees $T \in \mathcal{T}(s)$ to trees $T' \in \mathcal{T}(s')$ such that $\frac{\mathbf{P}(T)}{\mathbf{P}(T')} < \frac{\pi_{adj}(s)}{\pi_{adj}(s')}$. The basic mapping consists of reversing the path from s' to s in T to obtain the tree T' :

► **Definition 17.** For any s and s' , and for any tree $T \in \mathcal{T}(s)$ its s' -reversed is the tree T' equal to T but with the edges on the path from s' to s reversed.

Hereafter, we call such a tree T' simply a *reversed* tree of T if s and s' are clear from the context. We say that an edge has *length* w if its transition probability is of the form $\frac{\lambda^w}{\lambda^w + \lambda^{-w}}$. The *length* of a path is the total length of its edges. The *multiset* of a tree T is the set of absolute edge lengths appearing in T ,

$$ms(T) = \{w \mid \text{some edge in } T \text{ has length } w \text{ or } -w\} .$$

Based on this multiset, we denote by $ms(T, w)$ the number of edges in T whose length is either w or $-w$.

► **Fact 1.** Let T be a tree containing a path from s' to s of length ℓ , and let T' be its reversed tree. Then the probabilities of these two trees are of the form

$$\mathbf{P}(T) = \frac{\lambda^\ell \cdot \lambda^{L(T)}}{M(T)} \quad \text{and} \quad \mathbf{P}(T') = \frac{\lambda^{-\ell} \cdot \lambda^{L(T)}}{M(T)} \quad (14)$$

where $L(T)$ denotes the total length of the edges which are in the tree but not in the path between s and s' , and $M(T)$ depends on the multiset $ms(T)$, i.e.,

$$M(T) = \prod_{w \in ms(T)} (\lambda^w + \lambda^{-w})^{ms(T, w)} .$$

From (14) we get $\frac{\mathbf{P}(T)}{\mathbf{P}(T')} \leq \lambda^{2\ell}$. This motivates the following definition.

► **Definition 18.** For any s and s' , we call a tree $T \in \mathcal{T}(s)$ *good* if the length ℓ of its path from s' to s satisfies $\lambda^{2\ell} \leq \frac{\pi_{adj}(s)}{\pi_{adj}(s')}$. Otherwise we call T a *bad* tree.

By this definition the basic mapping of a bad tree does not give the desired bound (13). In such cases, we will map groups of bad trees into groups of good ones, depending on s and s' .

B.1.1.1 The (bac) vs (bca) case.

Observe that $\frac{\pi_{adj}(bac)}{\pi_{adj}(bca)} = \lambda^{2(c-a)}$. In this case there is no bad tree, since all paths from (bca) to (bac) have length at most $c - a$ (see Figure 7).

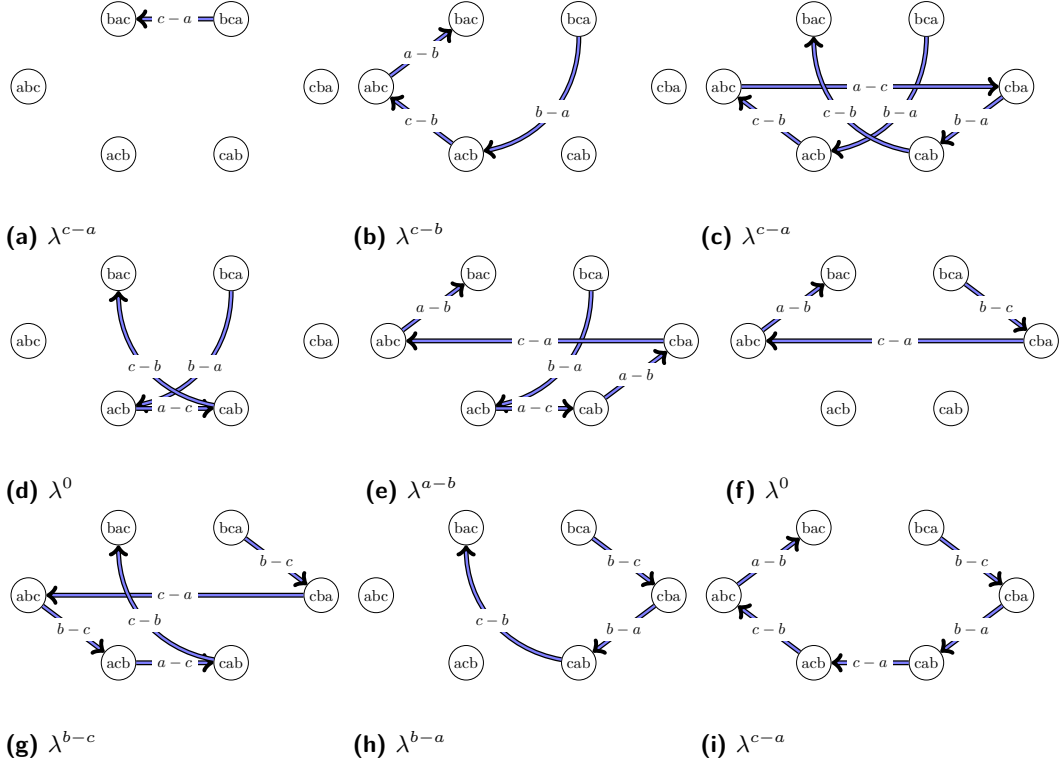
B.1.1.2 The (abc) vs (bac) case.

Note that $\frac{\pi_{adj}(bac)}{\pi_{adj}(bca)} = \lambda^{2(b-a)}$. Figure 8 shows all paths from $s' = (bac)$ to $s = (abc)$. Trees using the path in 8d are bad, since this path has length $\ell = c - a > b - a$. Trees using the path in 8c are bad if $c - b > b - a$. We combine the bad trees with the good trees having paths in 8b, 8f, 8h, or 8i.

► **Lemma 19.** Let $\mathcal{BAD}(s) \subset \mathcal{T}(s)$ be the set of bad trees with s' -path 8c or 8d, and let $\mathcal{GOOD}(s) \subset \mathcal{T}(s)$ be the set of good trees with 8b, 8f, 8h or 8i. Then,

$$\frac{\sum_{T \in \mathcal{BAD}(s)} \mathbf{P}(T) + \sum_{T \in \mathcal{GOOD}(s)} \mathbf{P}(T)}{\sum_{T \in \mathcal{BAD}(s)} \mathbf{P}(T') + \sum_{T \in \mathcal{GOOD}(s)} \mathbf{P}(T')} \leq \lambda^{2(b-a)} . \quad (15)$$

XX:16 Sort well with energy-constrained comparisons



■ **Figure 7** All paths from (bca) to (bac). There are no bad trees for these two states.

Proof. Let us use $x := b - a$, $y := c - b$, and $x + y := c - a$. The trees and their probabilities are shown in Figures 9-13 (in Appendix C). We will show

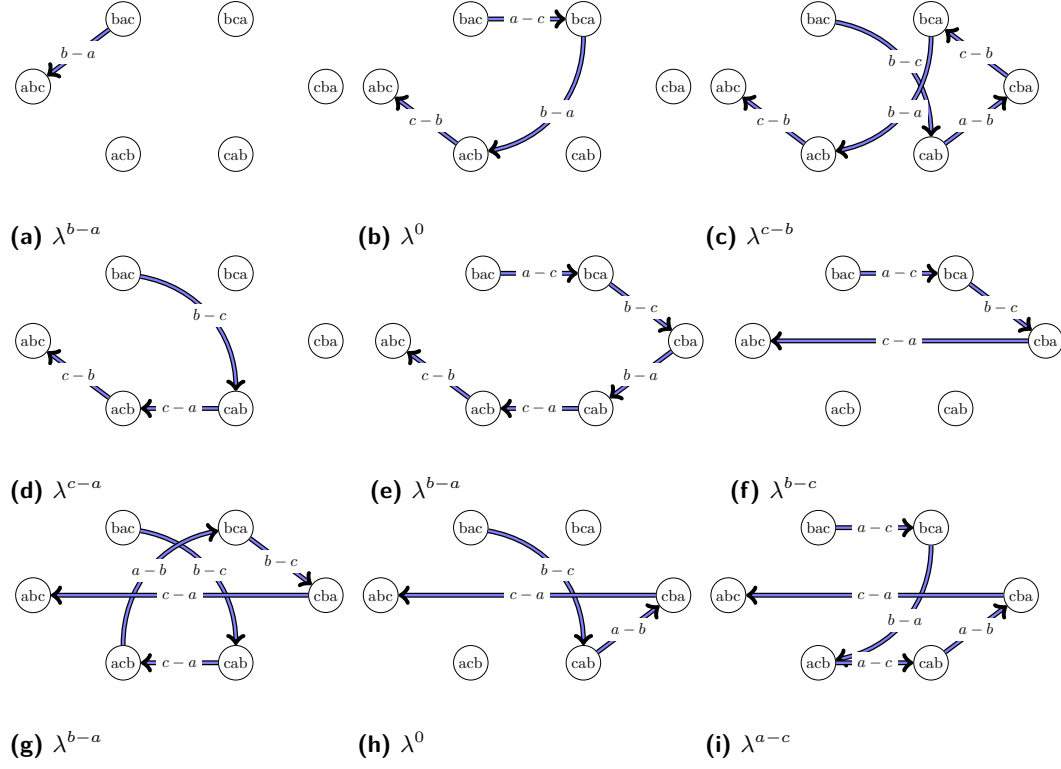
$$\left(\sum_{T \in \mathcal{BAD}(s)} \mathbf{P}(T) + \sum_{T \in \mathcal{GOOD}(s)} \mathbf{P}(T) \right) \leq \lambda^{2x} \left(\sum_{T \in \mathcal{BAD}(s)} \mathbf{P}(T') + \sum_{T \in \mathcal{GOOD}(s)} \mathbf{P}(T') \right)$$

which obviously implies (15). To get rid of the fractions in the probabilities we multiply them by the least common multiple of their denominators, i.e.,

$$\text{LCM} := (\lambda^{x+y} + \lambda^{-x-y})^3 (\lambda^y + \lambda^{-y})^3 (\lambda^x + \lambda^{-x})^2 .$$

We get for the left hand side of the inequality

$$\begin{aligned} \text{LCM} \left(\sum_{T \in \mathcal{BAD}(s)} \mathbf{P}(T) + \sum_{T \in \mathcal{GOOD}(s)} \mathbf{P}(T) \right) &= 6\lambda^{5x+4y} + 16\lambda^{3x+4y} + 10\lambda^{x+4y} + 6\lambda^{5x+2y} \\ &\quad + 26\lambda^{3x+2y} + 42\lambda^{x+2y} + 18\lambda^{-x+2y} + 8\lambda^{x-2y} \\ &\quad + 24\lambda^{-x-2y} + 12\lambda^{-3x-2y} + 4\lambda^{-x-4y} + 4\lambda^{-3x-4y} \\ &\quad + 10\lambda^{3x} + 36\lambda^x + 42\lambda^{-x} + 8\lambda^{-3x} . \end{aligned}$$



■ **Figure 8** All paths from (bac) to (abc). Bad trees include path (c) or (d).

And we get for the right hand side

$$\begin{aligned} \lambda^{2x} \text{LCM} \left(\sum_{T \in \text{BAD}(s)} \mathbf{P}(T') + \sum_{T \in \text{GOOD}(s)} \mathbf{P}(T') \right) = & 10\lambda^{5x+4y} + 16\lambda^{3x+4y} + 6\lambda^{x+4y} + 18\lambda^{5x+2y} \\ & + 42\lambda^{3x+2y} + 26\lambda^{x+2y} + 6\lambda^{-x+2y} + 8\lambda^{5x} \\ & + 42\lambda^{3x} + 36\lambda^x + 10\lambda^{-x} + 12\lambda^{3x-2y} \\ & + 24\lambda^{x-2y} + 8\lambda^{-x-2y} + 4\lambda^{x-4y} + 4\lambda^{-x-4y} . \end{aligned}$$

The difference between the left and the right hand side is

$$\begin{aligned} & 4\lambda^{5x+4y} - 4\lambda^{x+4y} + 12\lambda^{5x+2y} - 12\lambda^{-x+2y} + 16\lambda^{3x+2y} - 16\lambda^{x+2y} + 8\lambda^{5x} - 8\lambda^{-3x} \\ & + 32\lambda^{3x} - 32\lambda^{-x} + 12\lambda^{3x-2y} - 12\lambda^{-3x-2y} + 16\lambda^{x-2y} - 16\lambda^{-x-2y} + 4\lambda^{x-4y} - 4\lambda^{-3x-4y} \end{aligned}$$

and we can pair up the terms to conclude that this difference is positive. ◀

B.1.1.3 The (bca) vs (cba) case.

Note that $\frac{\pi_{adj}(bca)}{\pi_{adj}(cba)} = \lambda^{2(c-b)}$. Figure 14 in Appendix C shows all paths from $s' = (cba)$ to $s = (bca)$. We can combine the bad trees with s' - s -path 14d or 14e, and the good trees with paths 14c, 14f, 14h or 14i to show that the ratio between all trees and their reversals is smaller than $\lambda^{2(c-b)}$:

XX:18 Sort well with energy-constrained comparisons

► **Lemma 20.** *For the bad trees $\mathcal{BAD}(s) \subset \mathcal{T}(a)$ with s 's-path 14d or 14e, and the good trees $\mathcal{GOOD}(s) \subset \mathcal{T}(a)$ with paths 14c, 14f, 14h or 14i it holds that*

$$\frac{\sum_{T \in \mathcal{BAD}(s)} \mathbf{P}(T) + \sum_{T \in \mathcal{GOOD}(s)} \mathbf{P}(T)}{\sum_{T \in \mathcal{BAD}(s)} \mathbf{P}(T') + \sum_{T \in \mathcal{GOOD}(s)} \mathbf{P}(T')} \leq \lambda^{2(c-b)}.$$

Proof. We proceed as in the proof for Lemma 19 and finally get the difference between the left and the right hand side:

$$4\lambda^{4x+5y} - 4\lambda^{4x+y} + 12\lambda^{2x+5y} - 12\lambda^{2x-y} + 16\lambda^{2x+3y} - 16\lambda^{2x+y} + 8\lambda^{5y} - 8\lambda^{-3y} \\ + 32\lambda^{3y} - 32\lambda^{-y} + 12\lambda^{-2x+3y} - 12\lambda^{-2x-3y} + 16\lambda^{-2x+y} - 16\lambda^{-2x-y} + 4\lambda^{-4x+y} - 4\lambda^{-4x-3y}$$

◀

B.1.1.4 The other cases.

There are three other pairs of states for which we need to prove Lemma 13. However, the procedure is always the same: We identify the bad and good trees and combine them to conclude (10). For the missing paths consider Figures 15-17 in Appendix C.

B.2 One outlier

► **Lemma 15.** *The stationary distributions of \mathcal{M}_{adj} and \mathcal{M}_{any} , are*

$$\pi_{adj}(s^{(i)}) = \frac{p^{n-i}(1-p)^{i-1}(1-2p)}{(1-p)^n - p^n},$$

$$\pi_{any}(s^{(i)}) = \frac{np(1-p)}{((n-i+1)(1-p) + (i-1)p)((n-i)(1-p) + ip)}.$$

Proof. We get the stationary distribution for \mathcal{M}_{adj} using the global balance condition of stationary distribution¹, which implies that $\pi_{adj}(s^{(i)}) = \pi_{adj}(s^{(i+1)})(\frac{1-p}{p})$. By the structure of our Markov chain we get that for any state $s^{(j)}$ it holds that

$$\pi_{adj}(s^{(j)}) = \pi_{adj}(s^{(i)})(\frac{1-p}{p})^{i-j}.$$

Since the sum of all π_{adj} is one, we can express $\pi_{adj}(s^{(i)})$ as

$$\pi_{adj}(s^{(i)}) = 1 - \pi_{adj}(s^{(i)}) \left(\sum_{j=1}^{i-1} \left(\frac{p}{1-p}\right)^j + \sum_{j=1}^{n-i} \left(\frac{1-p}{p}\right)^j \right)$$

$$= \frac{1}{1 + \sum_{j=1}^{i-1} \left(\frac{p}{1-p}\right)^j + \sum_{j=1}^{n-i} \left(\frac{1-p}{p}\right)^j}.$$

By applying the formula for geometric series, we rewrite

$$\frac{1}{\pi_{adj}(s^{(i)})} = 1 + \frac{1 - (\frac{p}{1-p})^i}{1 - (\frac{p}{1-p})} - 1 + \frac{1 - (\frac{1-p}{p})^{n-i+1}}{1 - (\frac{1-p}{p})} - 1$$

$$= -1 + \frac{(1-p)^i - p^i}{(1-2p)(1-p)^{i-1}} + \frac{p^{n-i+1} - (1-p)^{n-i+1}}{(2p-1)p^{n-i}}$$

¹ The stationary distribution π of a Markov chain with transition matrix P must satisfy $\sum_{s' \neq s} \pi(s)P(s, s') = \sum_{s' \neq s} \pi(s')P(s', s)$, for any two states s, s' .

We expand all terms to the same denominator:

$$= \frac{-(1-2p)(1-p)^{i-1}p^{n-i} + (1-p)^i p^{n-i} - p^n - (1-p)^{i-1}p^{n-i+1} + (1-p)^n}{(1-2p)(1-p)^{i-1}p^{n-i}}$$

Finally, we can use that $(1-2p) = (1-p) - p$ and observe that the fraction simplifies to what we claimed.

The formula for the stationary distribution of \mathcal{M}_{any} can be derived as follows using the global balance condition of stationary distribution. For any $\pi_{any}(s^{(i)})$ it holds that

$$\begin{aligned} \pi_{any}(s^{(i)})((i-1)p + (n-i)(1-p)) &= (1-p) \sum_{j=1}^{i-1} \pi_{any}(s^{(j)}) + p \sum_{j=i+1}^n \pi_{any}(s^{(j)}) \\ &= (1-p) \sum_{j=1}^{i-1} \pi_{any}(s^{(j)}) + p \left(1 - \sum_{j=1}^i \pi_{any}(s^{(j)}) \right) \end{aligned}$$

that is

$$\pi_{any}(s^{(i)}) (ip + (n-i)(1-p)) = (1-2p) \sum_{j=1}^{i-1} \pi_{any}(s^{(j)}) + p .$$

Now we can show by induction on i that this recurrence resolves to

$$\pi_{any}(s^{(i)}) = \frac{np(1-p)}{((n-i+1)(1-p) + (i-1)p)((n-i)(1-p) + ip)} .$$

For $i = 1$ we immediately get $\pi_{any}(s^{(1)}) (p + (n-1)(1-p)) = p$, which we rewrite as

$$\pi_{any}(s^{(1)}) = \frac{p}{(n-1)(1-p) + p} = \frac{np(1-p)}{(n(1-p))((n-1)(1-p) + p)} .$$

For $i > 1$ we assume that the formula holds for $i-1$ and we get

$$\begin{aligned} \pi_{any}(s^{(i)}) &= \frac{(1-2p) \sum_{j=1}^{i-1} \pi_{any}(s^{(j)}) + p}{(ip + (n-i)(1-p))} \\ &= \frac{(1-2p) \sum_{j=1}^{i-2} \pi_{any}(s^{(j)}) + p + \pi_{any}(s^{(i-1)})(1-2p)}{(ip + (n-i)(1-p))} \\ &= \frac{\pi_{any}(s^{(i-1)})((i-1)p + (n-i+1)(1-p)) + \pi_{any}(s^{(i-1)})(1-2p)}{(ip + (n-i)(1-p))} \\ &= \frac{\pi_{any}(s^{(i-1)})((i-2)p + (n-i+2)(1-p))}{(ip + (n-i)(1-p))} \\ &= \frac{np(1-p)}{((n-i+1)(1-p) + (i-1)p)(ip + (n-i)(1-p))} . \end{aligned}$$

◀

XX:20 Sort well with energy-constrained comparisons

► **Lemma 16.** For \mathcal{M}_{adj} and \mathcal{M}_{any} , the corresponding expected total weighted inversions are

$$E_{adj}^w = n(b-a)p \left(\frac{1}{n(1-2p)} - \frac{p^{n-1}}{(1-p)^n - p^n} \right) < (b-a) \frac{p}{1-2p} ,$$

$$E_{any}^w = n(b-a)p \cdot \sum_{i=0}^{n-1} \frac{i(1-p)}{((i+1)(1-p) + (n-i-1)p)(i(1-p) + (n-i)p)} > n(b-a)p .$$

Proof. We apply the generic formula for the expected weighted inversion (12) to derive the expected weighted inversion of \mathcal{M}_{adj} and \mathcal{M}_{any} . Since $\pi_{adj}(s^{(i)}) = \pi_{adj}(s^{(n)}) \left(\frac{p}{1-p} \right)^{n-i}$ we get

$$\begin{aligned} E_{adj}^w &= \sum_{i=1}^n (n-i)(x-1) \pi_{adj}(s^{(i)}) \\ &= (x-1) \pi_{adj}(s^{(n)}) \cdot \sum_{i=0}^{n-1} i \left(\frac{p}{1-p} \right)^i \\ &= (x-1) \frac{(1-p)^{n-1}(1-2p)}{(1-p)^n - p^n} \cdot \frac{(n-1) \left(\frac{p}{1-p} \right)^{n-1} - n \left(\frac{p}{1-p} \right)^n + \left(\frac{p}{1-p} \right)}{\left(\frac{p}{1-p} - 1 \right)^2} \\ &= n(x-1)p \left(\frac{1}{n(1-2p)} - \frac{p^{n-1}}{(1-p)^n - p^n} \right) . \end{aligned}$$

Then observe that for $0 < p < \frac{1}{2}$, the first inequality is immediate from $\frac{p^n}{(1-p)^n - p^n} > 0$. As for \mathcal{M}_{any} we have

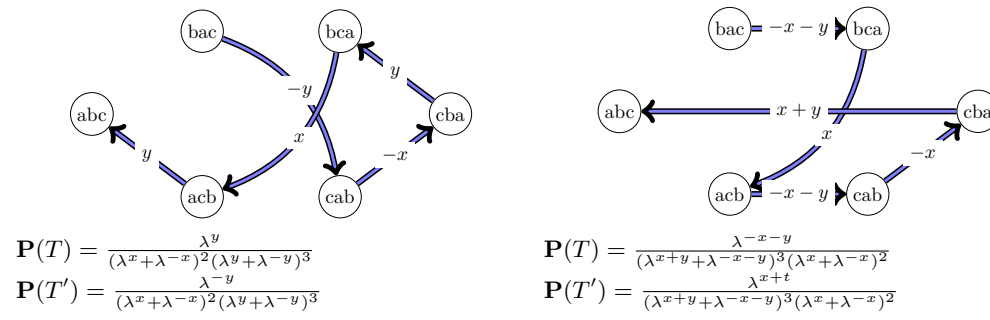
$$\begin{aligned} E_{any}^w &= \sum_{i=1}^n (n-i)(x-1) \pi_{any}(s^{(i)}) \\ &= (x-1) \cdot \sum_{i=0}^{n-1} i \cdot \pi_{any}(s^{(n-i)}) \\ &= n(x-1)p \cdot \sum_{i=0}^{n-1} \frac{i(1-p)}{((i+1)(1-p) + (n-i-1)p)(i(1-p) + (n-i)p)} . \end{aligned}$$

Observe that we can lower bound the sum in the formula by the integral

$$\int_0^n \frac{i(1-p)}{((i+1)(1-p) + (n-i-1)p)(i(1-p) + (n-i)p)} di ,$$

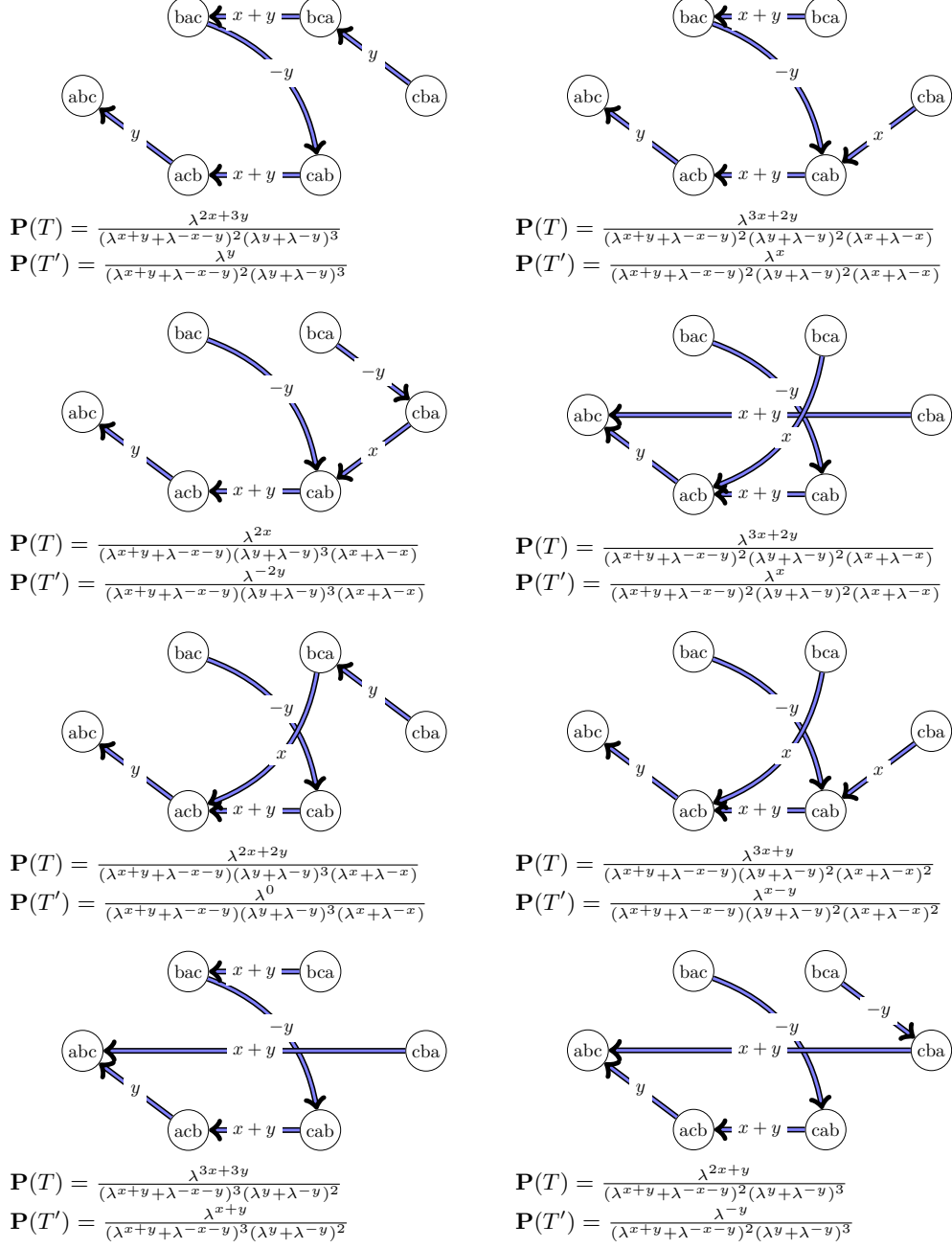
which is larger than 1 if $0 < p < \frac{1}{2}$. ◀

C Additional pictures

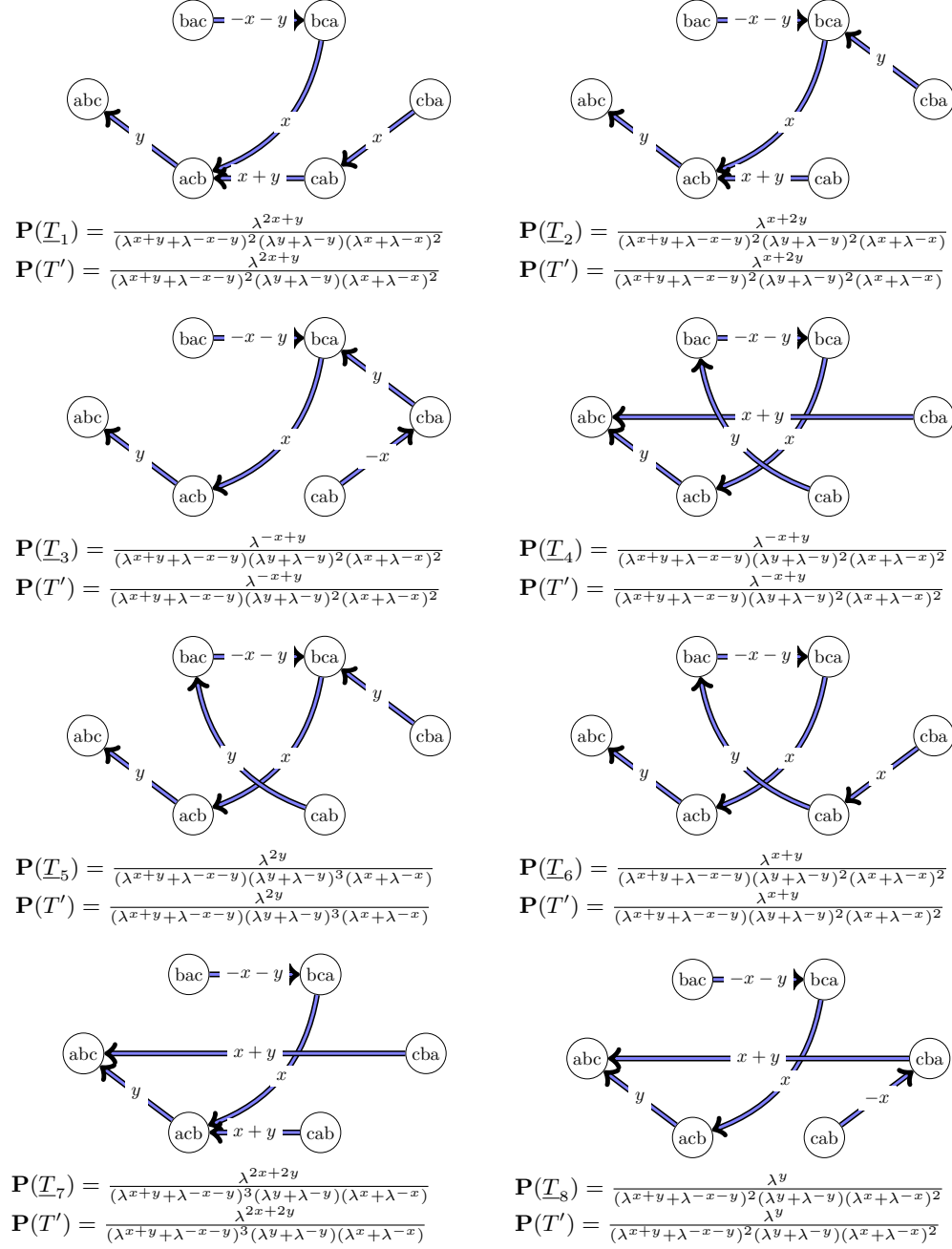


■ **Figure 9** Bad tree with path 8c and good tree with path 8i.

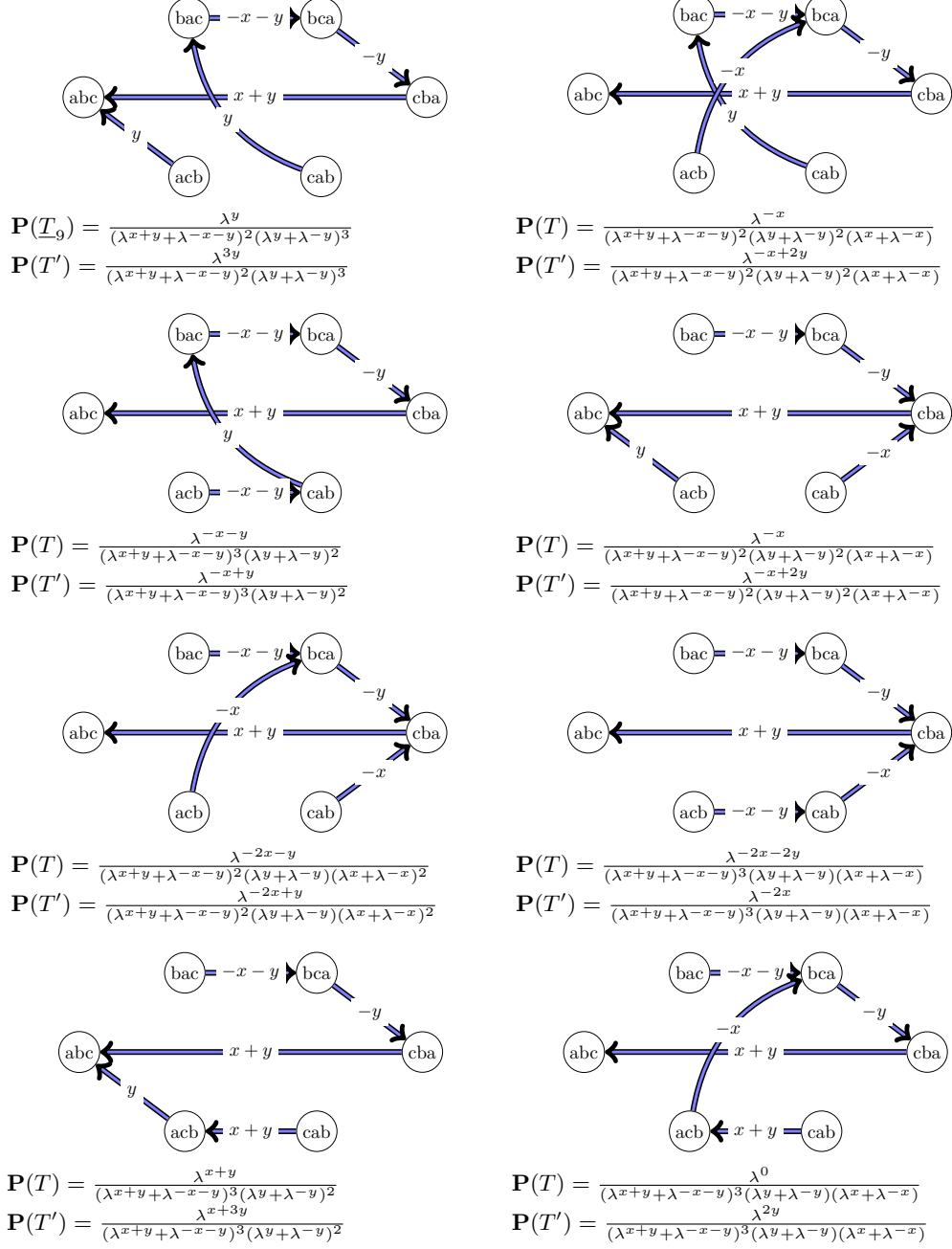
XX:22 Sort well with energy-constrained comparisons



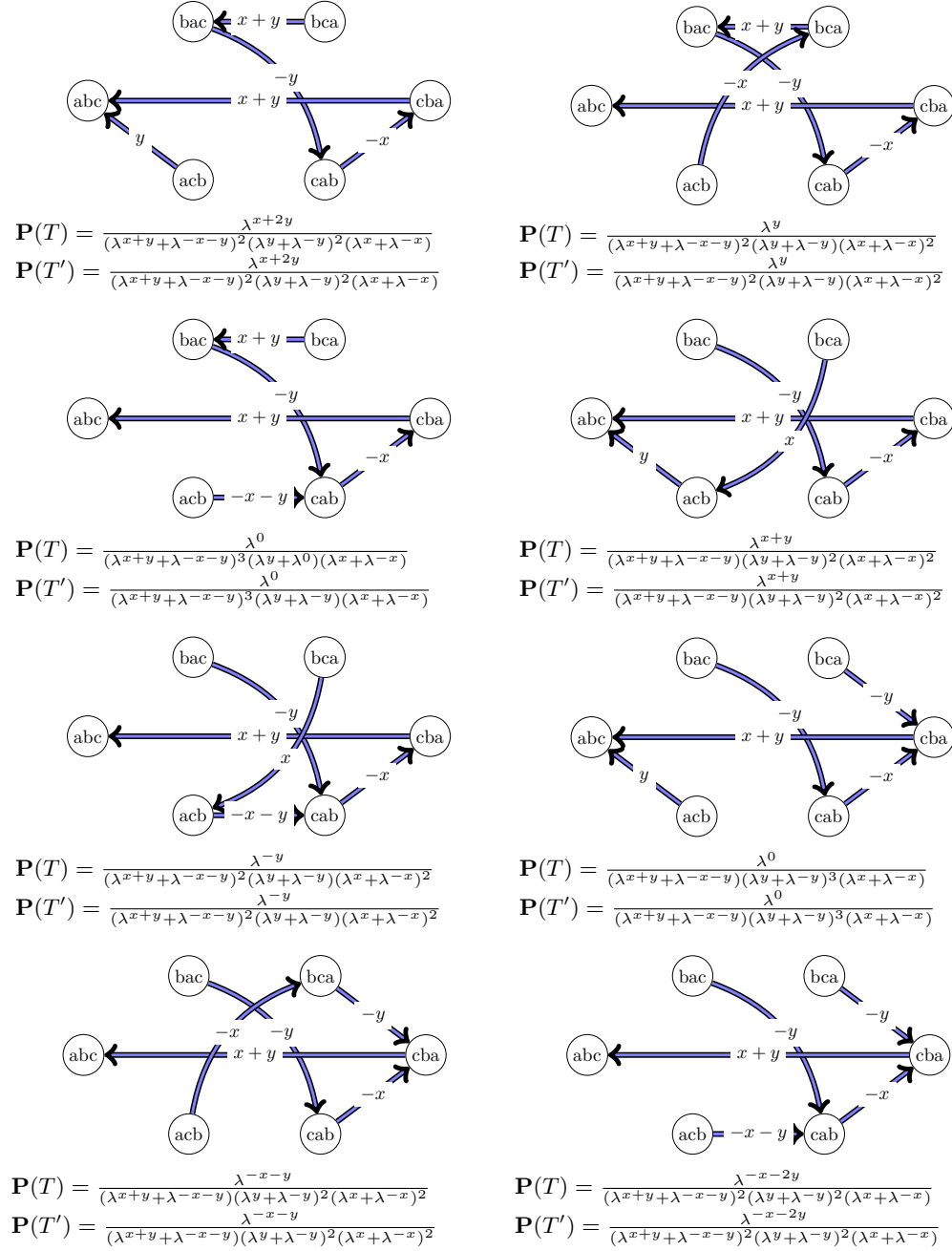
■ **Figure 10** Bad trees with path 8d.



■ **Figure 11** Good trees with path 8b.

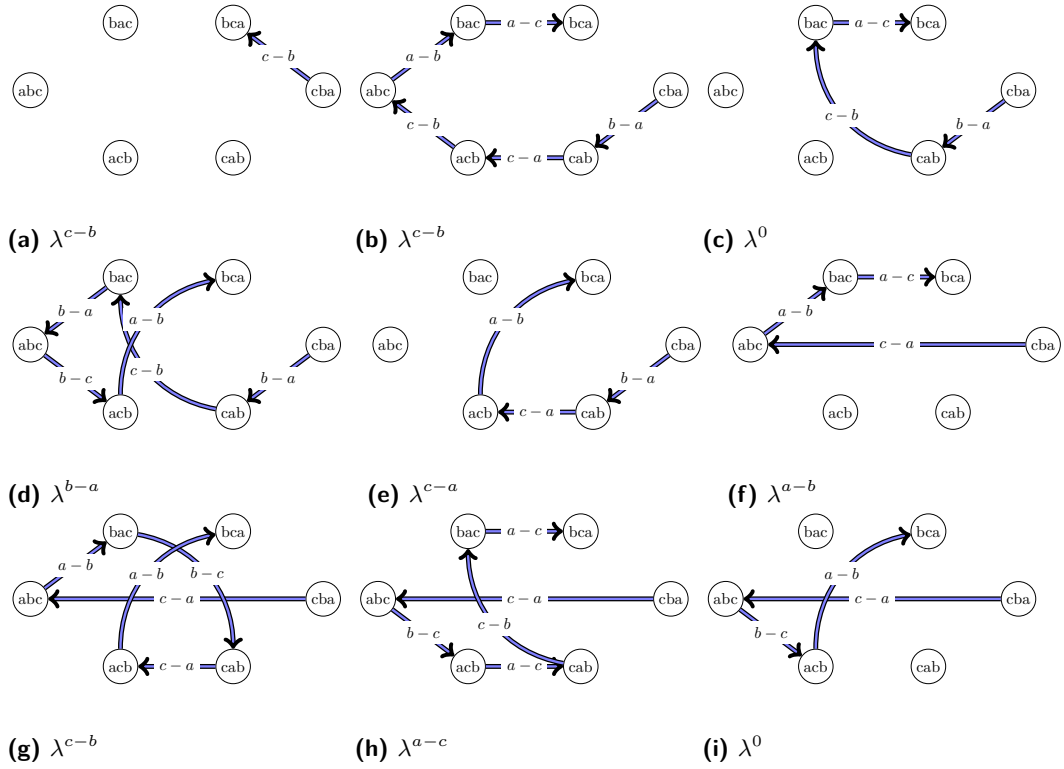


■ Figure 12 Good trees with path 8f.

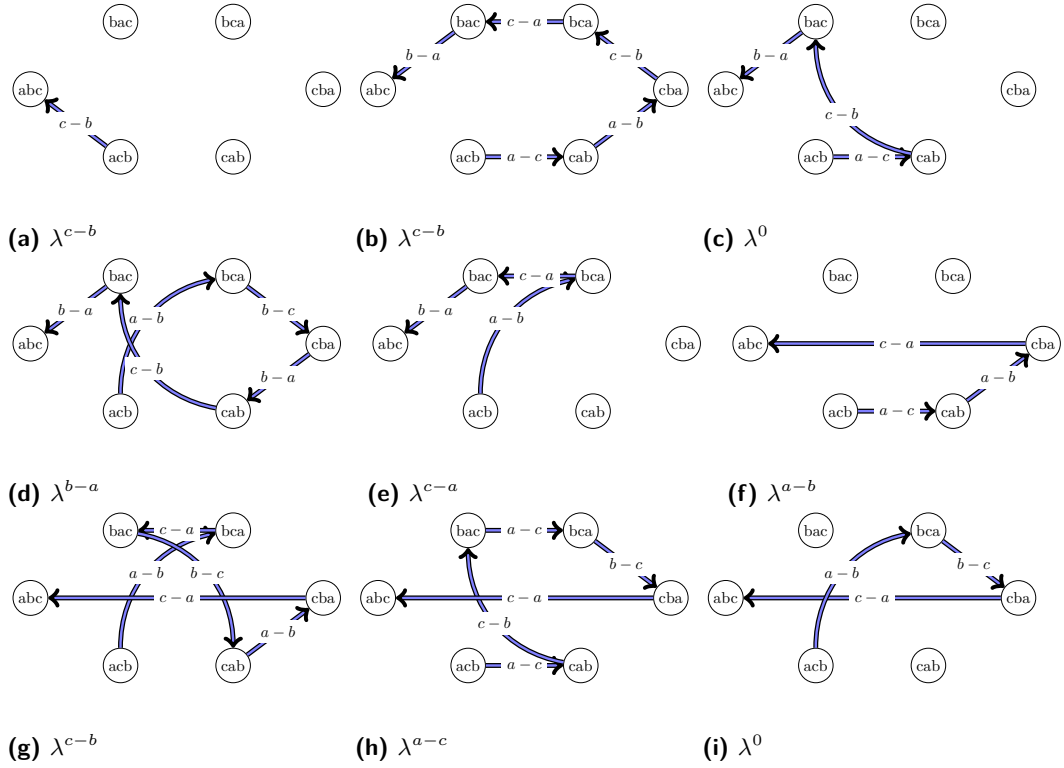


■ **Figure 13** Good trees with path 8h.

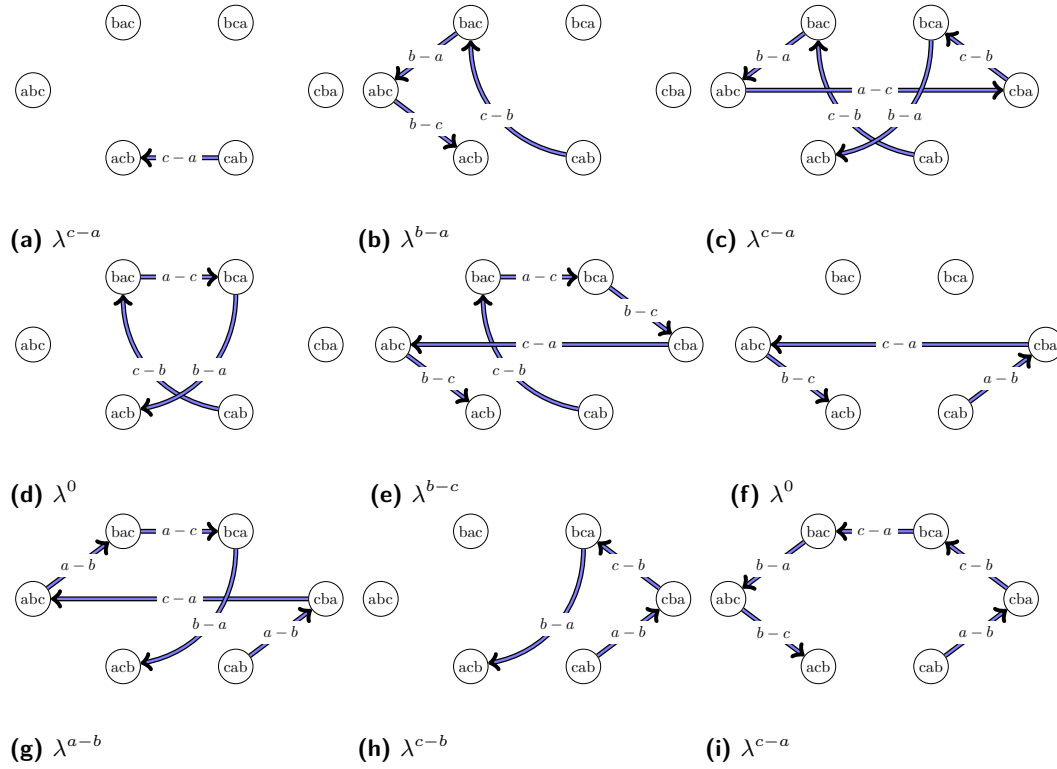
XX:26 Sort well with energy-constrained comparisons



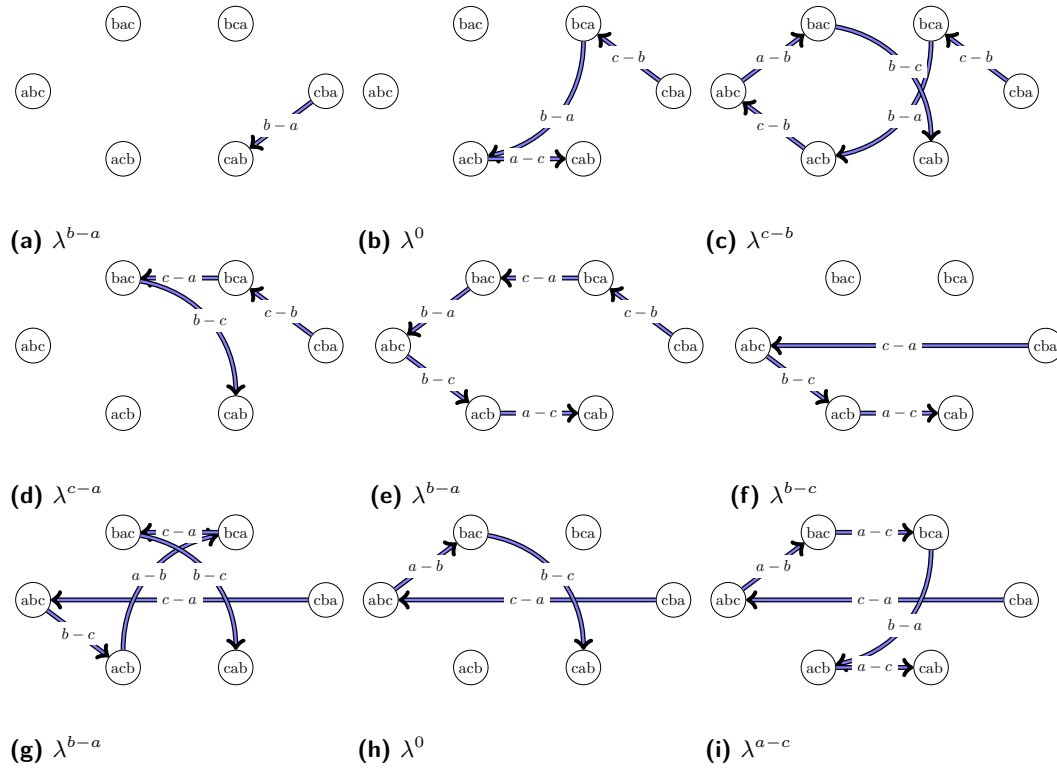
■ **Figure 14** All paths from (cba) to (bca). Bad trees include path (d) or (e).



■ **Figure 15** All paths from (acb) to (abc). Bad trees include path (d) or (e).



■ **Figure 16** All paths from (cab) to (acb). There are no bad trees for these two states.



■ **Figure 17** All paths from (cba) to (cab). Bad trees include path (c) or (d).